

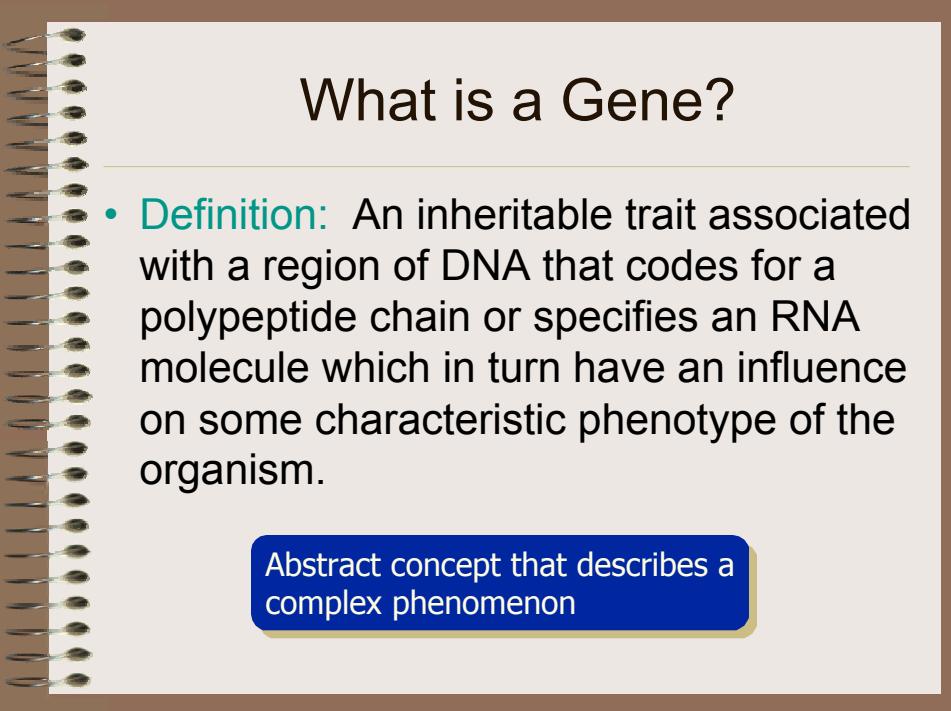
Gene finding and Genome annotation

Manfred Zorn

BerkeleyPGA

Bioinformatics Tools for Comparative Analysis

June 12, 2001



What is a Gene?

- **Definition:** An inheritable trait associated with a region of DNA that codes for a polypeptide chain or specifies an RNA molecule which in turn have an influence on some characteristic phenotype of the organism.

Abstract concept that describes a complex phenomenon

What is Annotation?

- **Definition:** Extraction, definition, and interpretation of features on the genome sequence derived by integrating computational tools and biological knowledge.

Identifiable features in the sequence

How does an annotation differ from a gene?

- Many annotations describe features that constitute a gene.
- Other annotations may not always directly correspond in this way, e.g., an STS, or sequence overlap



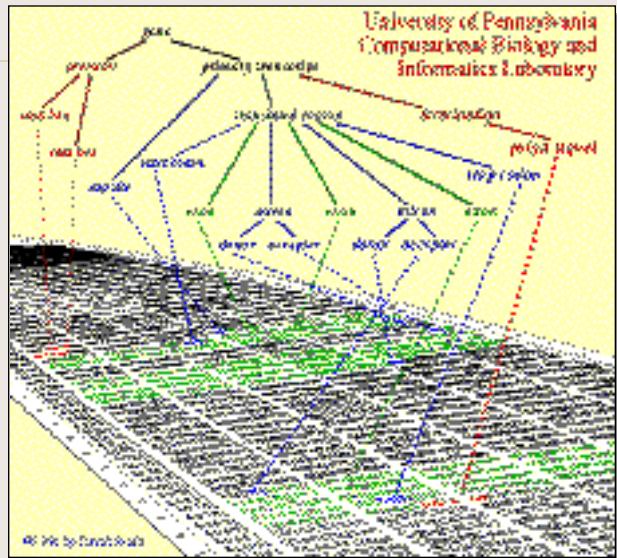
DNA Analysis

- Heuristics
- Statistics
- Artistics

DNA Analysis

- Find the genes
 - Heuristic signals
 - Inherent features
 - Intelligent methods
- Characterize each gene
 - Compare with other genes
 - Find functional components
 - Predict features

What is a Gene?



Heuristic Signals

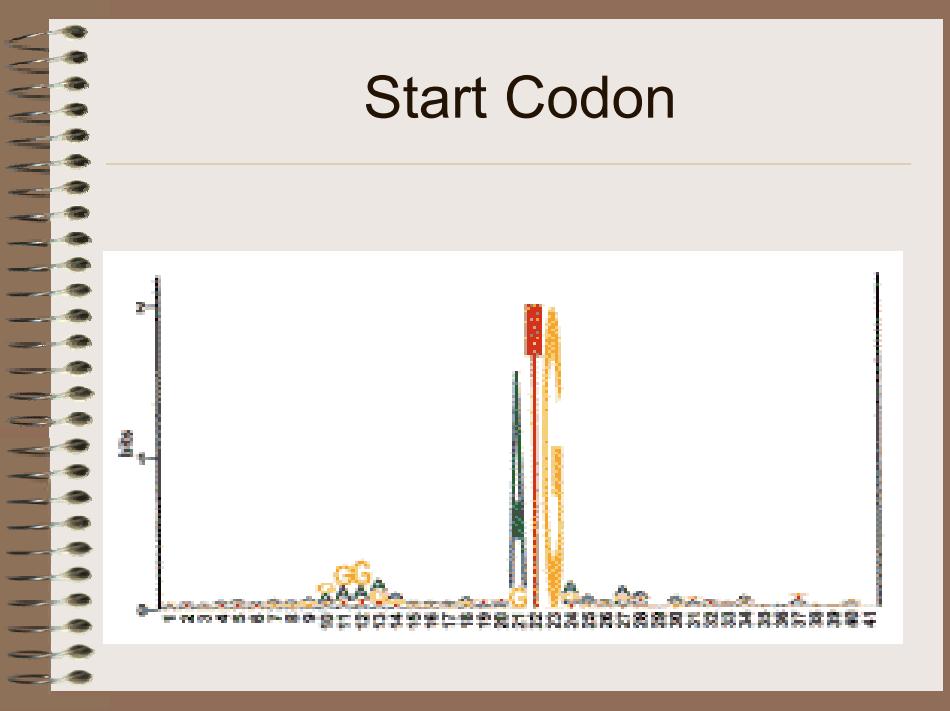
- DNA contains various recognition sites for internal machinery
- Promoter signals
- Transcription start signals
- Start Codon
- Exon, Intron boundaries
- Transcription termination signals

Heuristic Signals

Start of the gene

Heuristic Signals

Start Codon



Initiation

AGCACBGGGAAAUCUGAUGAACGUAC *E. coli* rplB
UUUGGAGGGGUGAAACGAGUGCAUUGCA *E. coli* rplB
GGUAACGGGAAACAACCAGAGCGAGUGU *E. coli* rplB
CAAUCGGGGUGGAAUGUAACCAGUA *E. coli* rplB
AAUCUGGGABGCUUUUUAUGUCUUU *φX174* phage A protein
UAACUGGGAAGCAGUCUAAGACA *Q*β phage replicase
UCCUGGGGGGURUGACCAUGCGGCUUU *R*17 phage A protein
AUGUACGGGGGGGGGUUUGUAAACAACGC *λ* phage cII

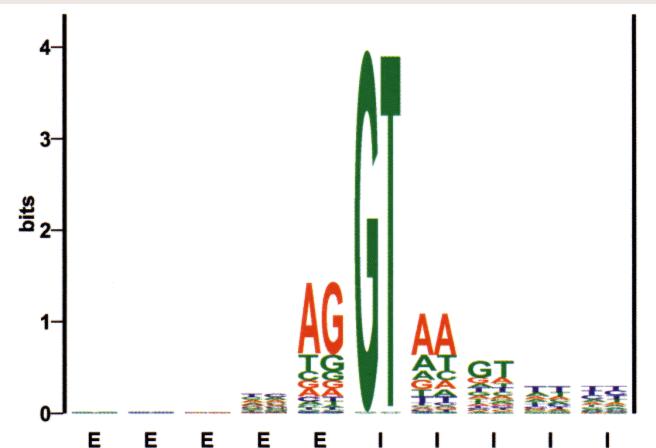
Paired with 16S rRNA Paired with 23S rRNA



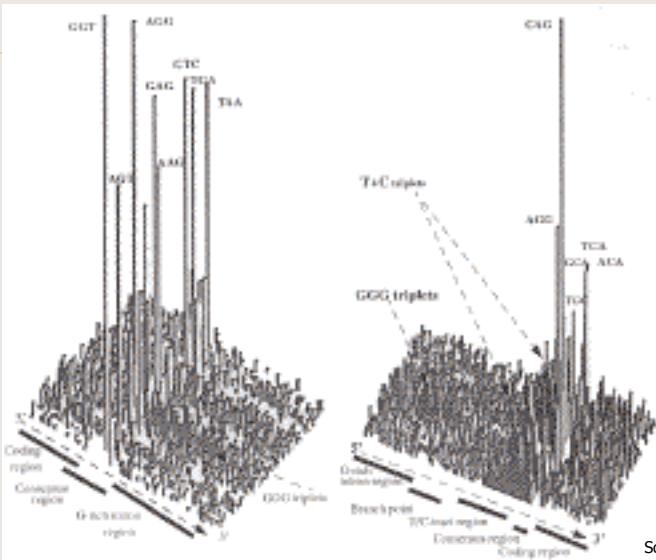
Inherent Features

- DNA exhibits certain biases that can be exploited to locate coding regions
- Uneven distribution of bases
- Codon bias
- CpG islands
- In-phase words
- Encoded amino acid sequence
- Imperfect periodicity
- Other global patterns

Donor Splice Site



Inherent Features



Intelligent Methods

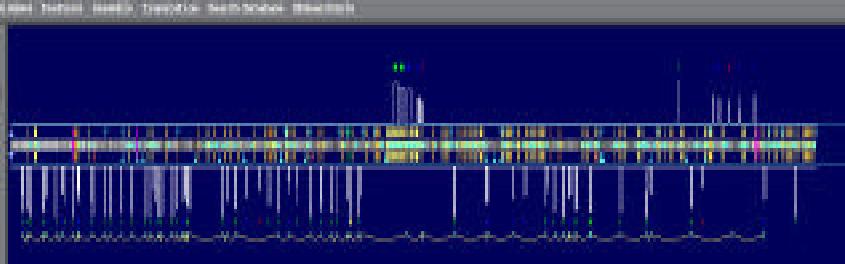
- Pattern recognition methods weigh inputs and predict gene location
 - Content-based methods
 - Site-based methods
 - Comparative methods
- Neural Networks
- Hidden Markov Models
- Stochastic Context-Free Grammar

GRAIL *Uberbacher, Mural*

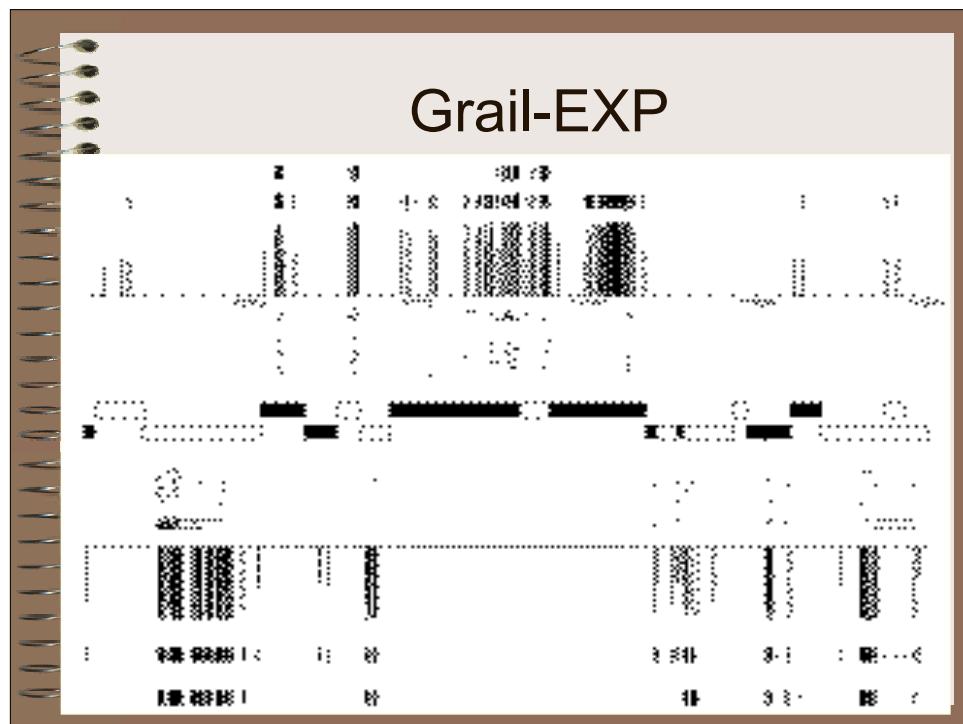
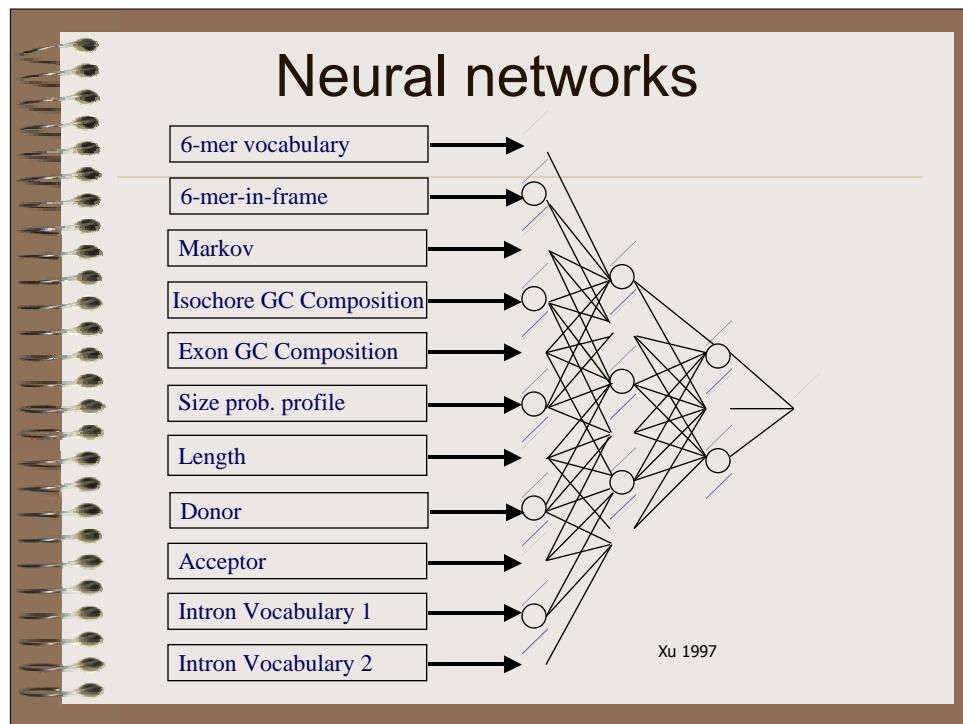
- GRAIL 1
 - Neural network with fixed window length (100 bases)
- GRAIL 1a
 - GRAIL 1 + adjacent information
- GRAIL 2
 - Variable length window, contextual information
- GRAIL-EXP
 - Comparison with partial and complete gene sequences

Analyzing Complex Multi-Gene Regions

- Errors in exon prediction and splice site boundaries
- Gene boundaries uncertain
- Genes can be on both strands



Überbacher



FGNEH/FGENES Solovyev

- Looks at several structural features
 - Splice donor/acceptor sites
 - Putative coding regions
 - Intronic regions
- *Linear discriminant analysis* to split exon / non-exon classes
- Dynamic programming to assemble best gene structure

MZEF Zhang

- *Quadratic discriminant analysis*
 - Exon length
 - Exon-intron transitions
 - Splice sites
 - Branch sites
 - Exon, strand, frame scores
- Detects internal exons
- No information about gene structure



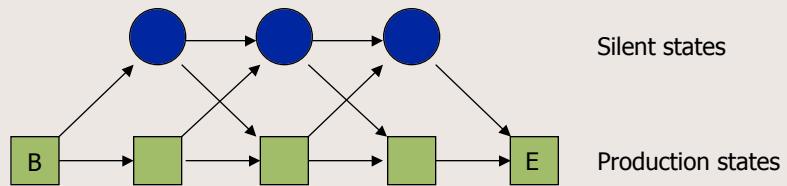
GENSCAN *Burge, Karlin*

- Probabilistic model of sequence composition and gene structure
 - Looks for gene structure descriptions that are consistent with the query sequence to assign probability that sequence stretch is exon, ...
 - Best ---> optimal
 - But generates also suboptimal exons

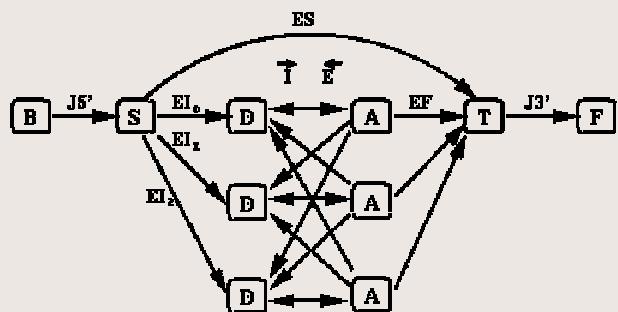
PROCRUSTES *Gelfand*

- Forces sequence into target structure
 - Requires putative gene product
 - Stretches/shortens sequence to fit into model

Hidden Markov Models



GENIE *Kulp, Reese, Haussler*





Strategies

-
- Select by correlation coefficient
 - Select by review paper
 - Select by recommendation
 - Use them all



Drawbacks

-
- Most programs are “trained” on existing data
 - It’s awfully hard to find new things this way!
 - NTT
 - IPW

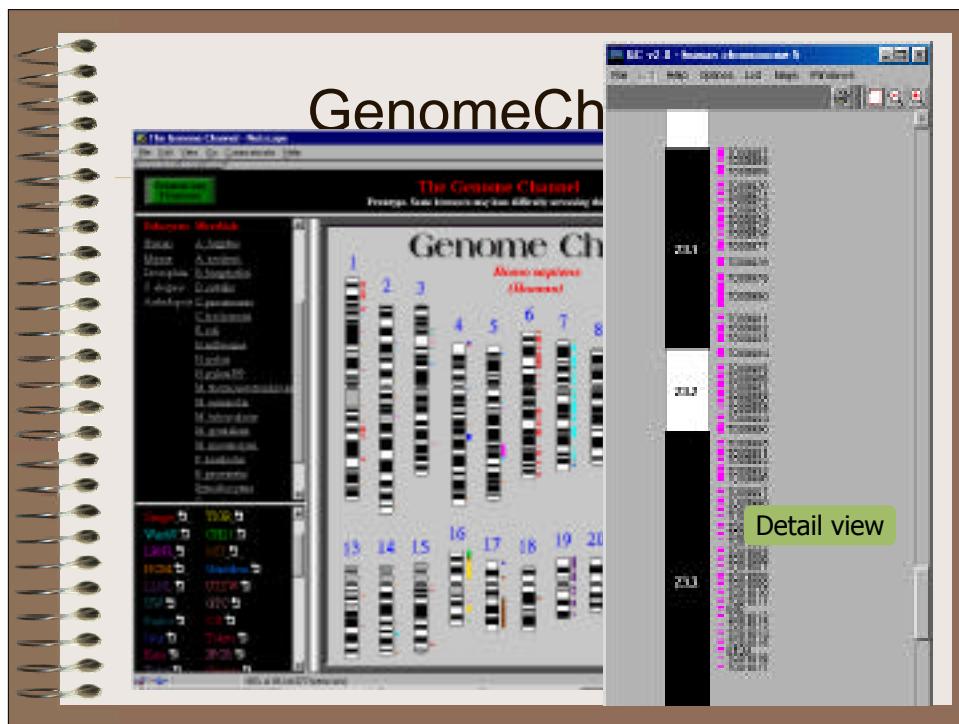
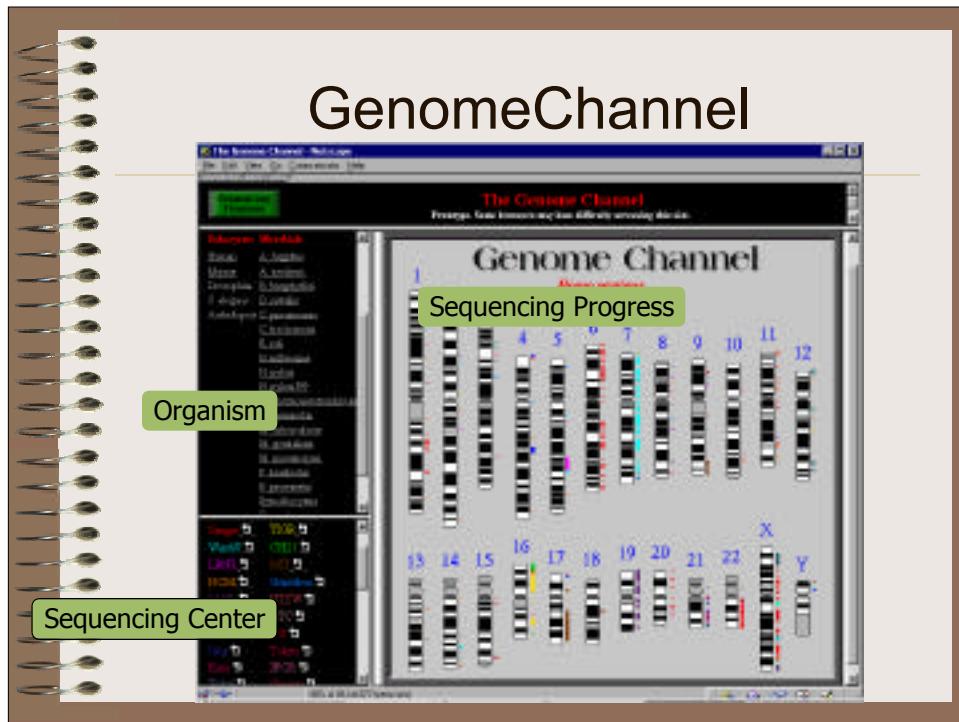
Internet Resources

| | |
|---------------|---|
| Banbury Cross | http://igs-server.cnrs-mrs.fr/igs/banbury |
| FGENEH | http://genomic.sanger.ac.uk/gf/gf.shtml |
| GenID | http://www1.imim.es/geneid.html |
| GeneMachine | http://genome.nhgri.nih.gov/genemachine |
| GENSCAN | http://genes.mit.edu/GENSCAN.html |
| Genotator | http://www.fruitfly.org/_nomi/genotator/ |
| GRAIL | http://compbio.ornl.gov/tools/index.shtml |
| GRAIL-EXP | http://compbio.ornl.gov/grailexp |
| MZEF | http://www.cshl.org/genefinder |
| PROCRUSTES | http://www-hto.usc.edu/software/procrustes |
| RepeatMasker | http://ftp.genome.washington.edu/RM/RepeatMasker.html |
| HMMgene | http://www.cbs.dtu.dk/services/HMMgene |
| Chapter 10 | http://www.wiley.com/legacy/products/subject/life/bioinformatics/chapterlinks.html |

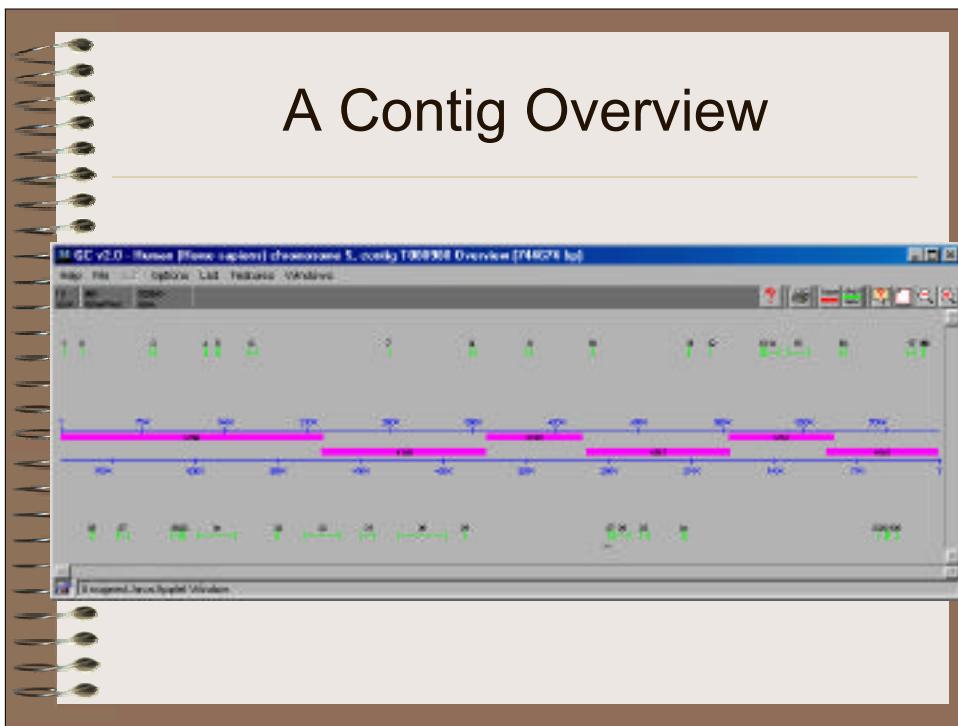
Characterize a Gene

Collect clues for potential function

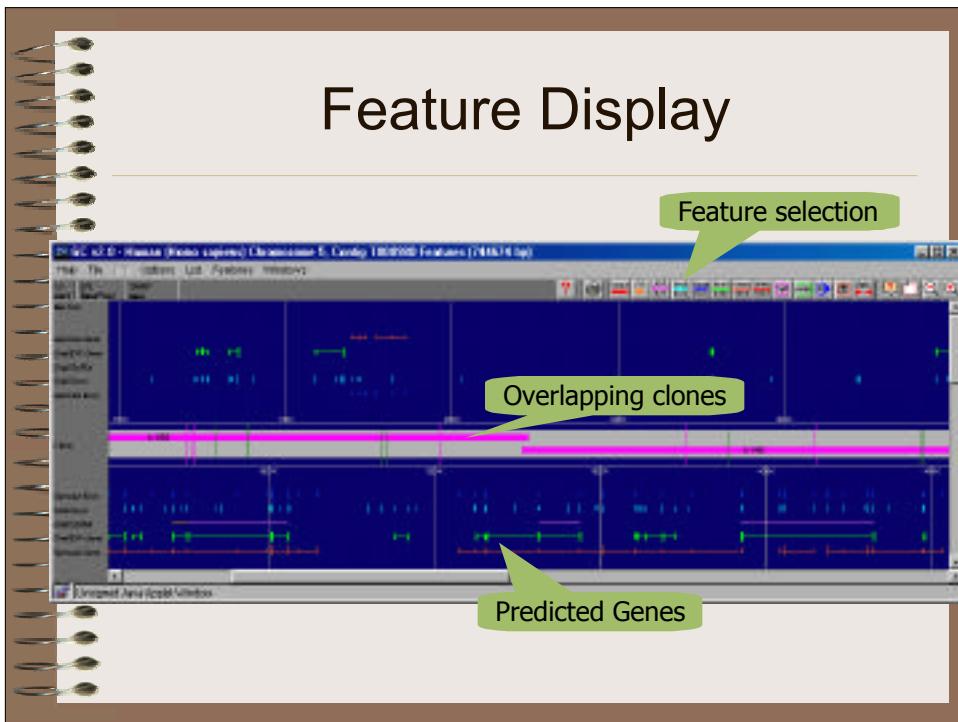
- Comparison with other known genes, proteins
- Predict secondary structure
- Fold classification
- Gene Expression
- Gene Regulatory Networks
- Phylogenetic comparisons
- Metabolic pathways

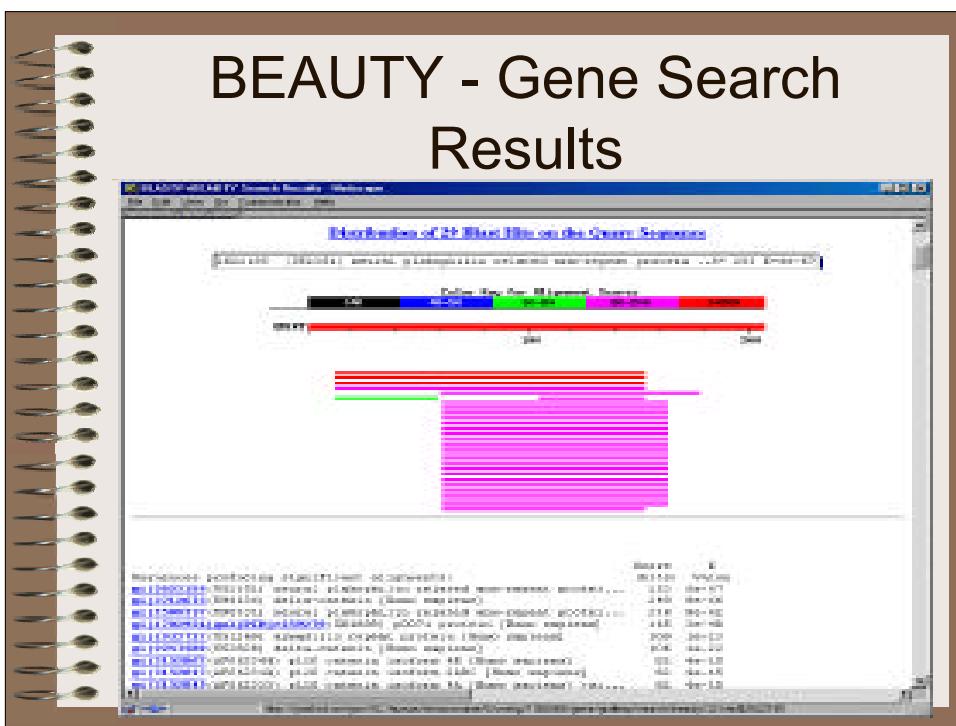
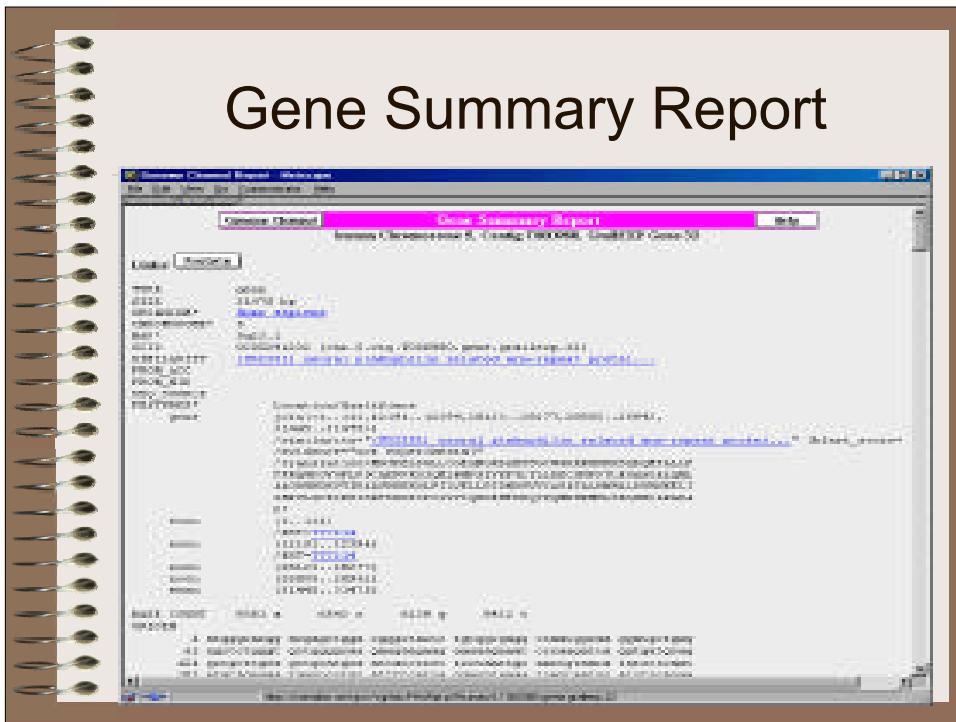


A Contig Overview



Feature Display



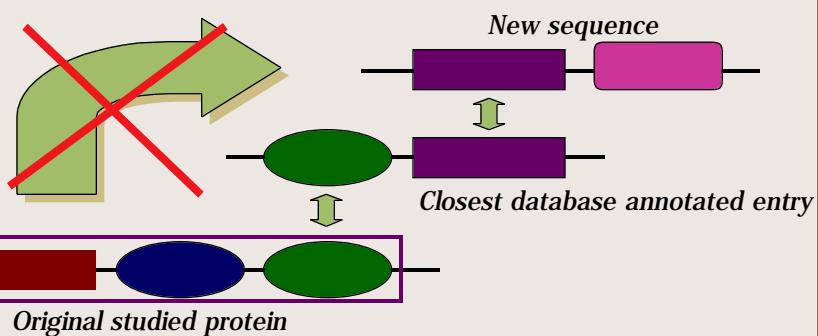


Layers of Information

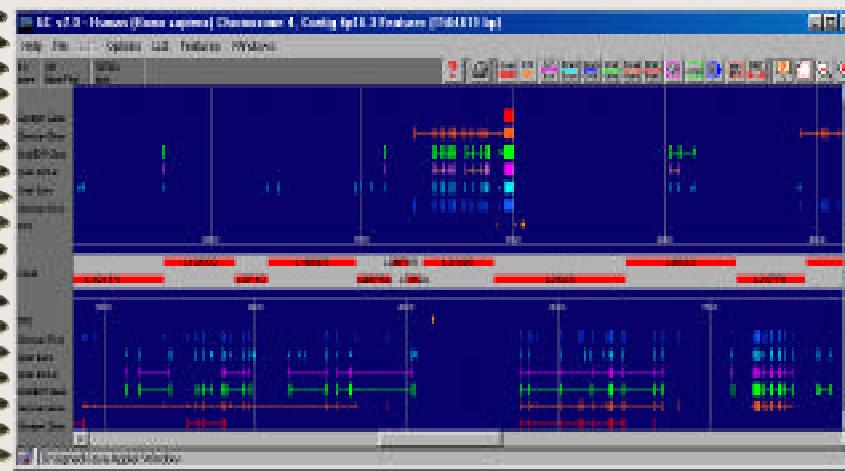
The same base sequence contains many layered instructions!

- Chromosome structure and function
 - Telomers, centromers
- Gene Regulatory information
 - Enhancers, promoters, ...
- Instructions for gene structure
- Instructions for protein
- Instructions for protein post-processing and localization

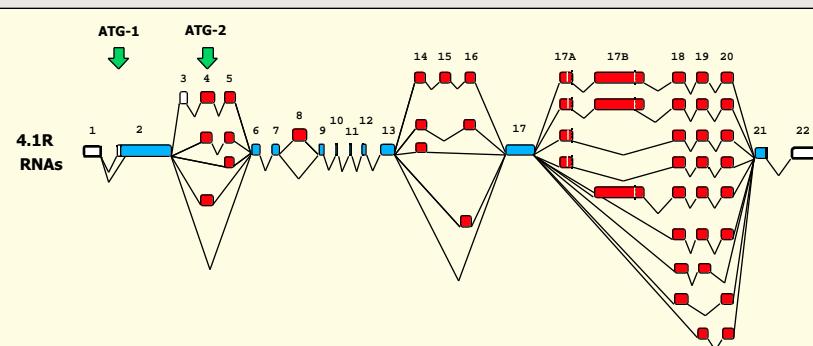
Inherited Annotation Problems in Multi-Domain Proteins



Alternatively Spliced ?

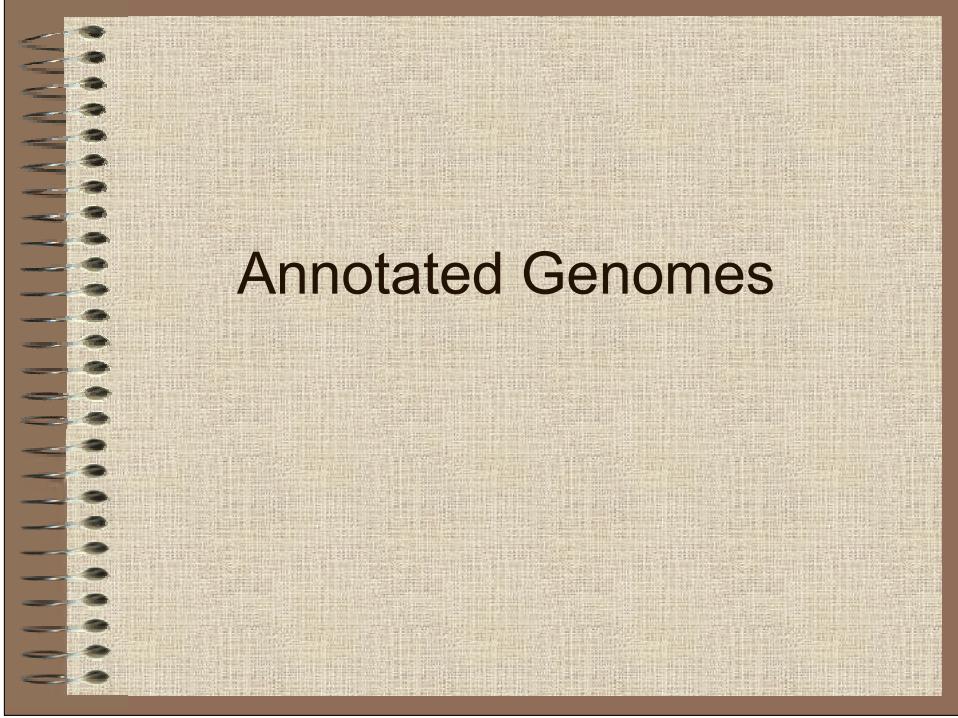


One Gene - Many Proteins



Conboy 1998

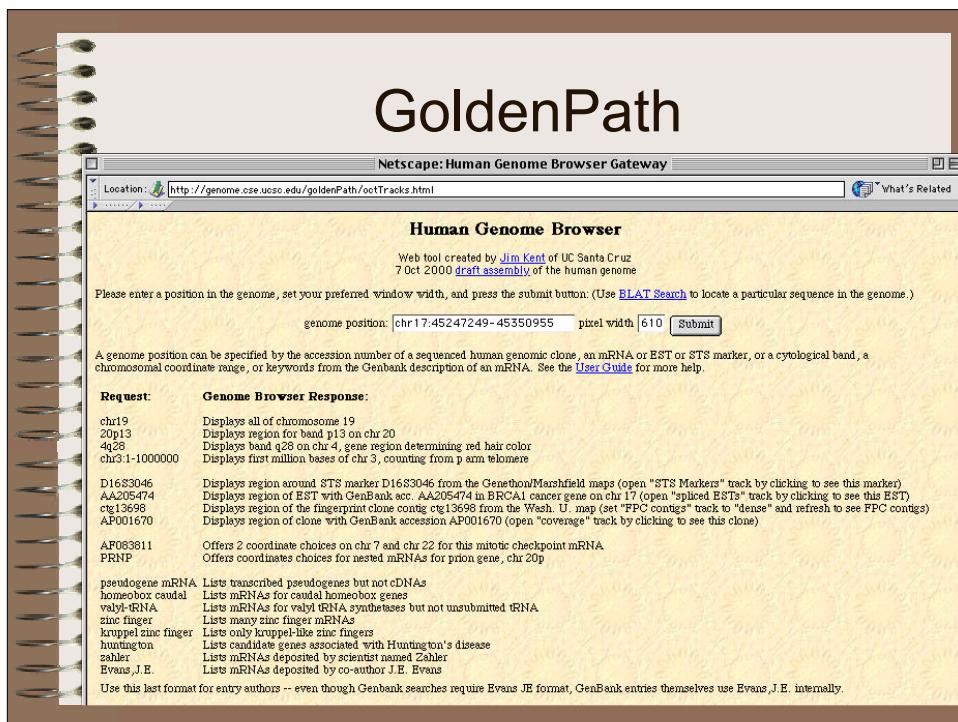
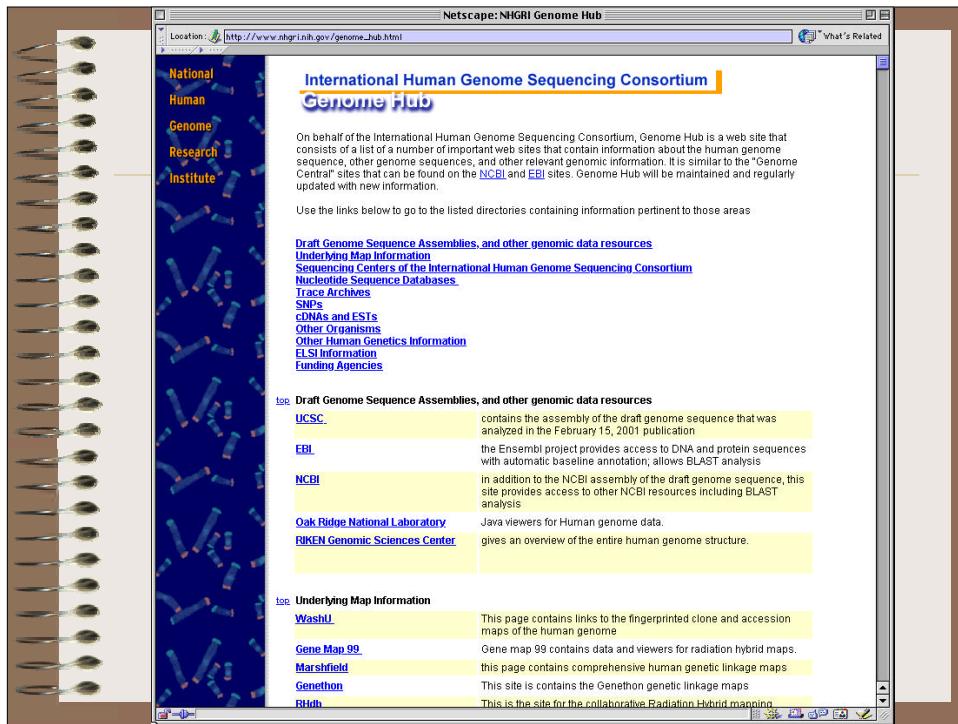
As many as 30% of human genes,
in particular structural genes, may
be alternatively spliced.



Annotated Genomes

Internet Resources

-
- | | |
|--------------------|---|
| GenomeHub | http://www.nhgri.nih.gov/genome_hub.html |
| NAR Database Issue | http://nar.oupjournals.org/content/vol29/issue1/ |





BLAT Search

Netscape: BLAT Search

Location: http://genome.cse.ucsc.edu/cgi-bin/hgBlat?db=hg5 What's Related

BLAT Search Human Genome

Freeze: Oct 7, 2000 Query type: BLAT's guess Sort output: query_score Submit

Please paste in a query sequence to see where it is located in the UCSC assembly of the human genome. Multiple sequences can be searched at once if separated by a line starting with > and the sequence name.

Only DNA sequences less than 20,000 bases and protein or translated sequence of less than 4000 letters will be processed. If multiple sequences are submitted at the same time, the total limit is 50,000 bases or 10,000 letters.

BLAT on DNA is designed to quickly find sequences of 95% and greater similarity of length 40 bases or more. It may miss more divergent or shorter sequence alignments. It will find perfect sequence matches of 33 bases, and sometimes find them down to 22 bases. BLAT on proteins finds sequences of 80% and greater similarity of length 20 amino acids or more. In practice DNA BLAT works well on primates, and protein blot on land vertebrates.

BLAT is not BLAST. DNA BLAT works by keeping an index of the entire genome in memory. The index consists of all non-overlapping 11-mers except for those heavily involved in repeats. The index takes up a bit less than a gigabyte of RAM. The genome itself is not kept in memory, allowing BLAT to deliver high performance on a reasonably priced Linux box. The index is used to find areas of probable homology, which are then loaded into memory for a detailed alignment. Protein BLAT works in a similar manner, except with 4-mers rather than 11-mers. The protein index takes a little more than 2 gigabytes.

BLAT was written by Jim Kent. Like most of Jim's software interactive use on the web server is free to all. Sources and executables to run batch jobs on your own server are available free for academic, personal, and non-profit purposes. Non-exclusive commercial licenses are also available. Contact Jim for details.

NCBI Home > Genomic Biology > Human

Search LocusLink for Go

The Human Genome

A guide to online information resources

Web Resources

- BLAST.** Compare your sequence to the genome or its gene products.
- Cytogenetics.** A cytogenetic resource of FISH-mapped, sequence-tagged clones.
- dbSNP.** Database of SNPs and other genetic variations.
- e-PCR.** Check your sequence for STSs and view in genomic context.
- GEO.** Gene Expression Omnibus, a public repository for expression data.
- HomoloGene.** Putative homologies among human, mouse, rat, and zebrafish.
- Homology Map.** Blocks of conserved synteny between mouse and human.
- LocusLink.** Focal point for genes and associated information.
- OMIM.** Guide to genes and inherited disorders maintained by JHU and collaborators.
- RefSeq.** Reference sequences of

Building an information infrastructure

A challenge facing researchers today is the ability to piece together and analyze the multitudes of data currently being generated through the Human Genome Project. NCBI's Web site serves as an integrated, one-stop, genomic information infrastructure for biomedical researchers from around the world so that they may use this data in their research efforts. [More...](#)

Working Draft Analysis Published

- NLM Press Release
- NHGRI Press Release
- Interactive Tour of the Genome
- NCBI Genome Analysis Pipeline
- Nature (2/15/01) Human Genome Issue
- Science (2/16/01) Human Genome Issue

Browse

Genes

1 2 3 4 5 6 7 8
9 10 11 12 13 14 15 16
17 18 19 20 21 22 X Y

Genes & Disease

G&D. Selected gene stories for students and the public.

MapViewer tips and tricks

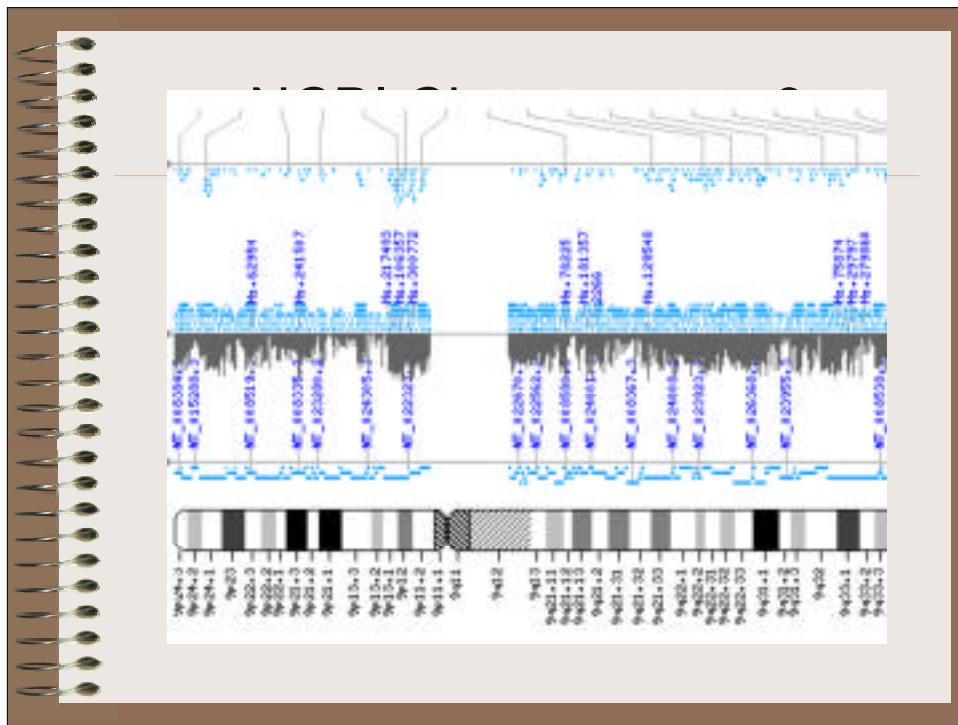
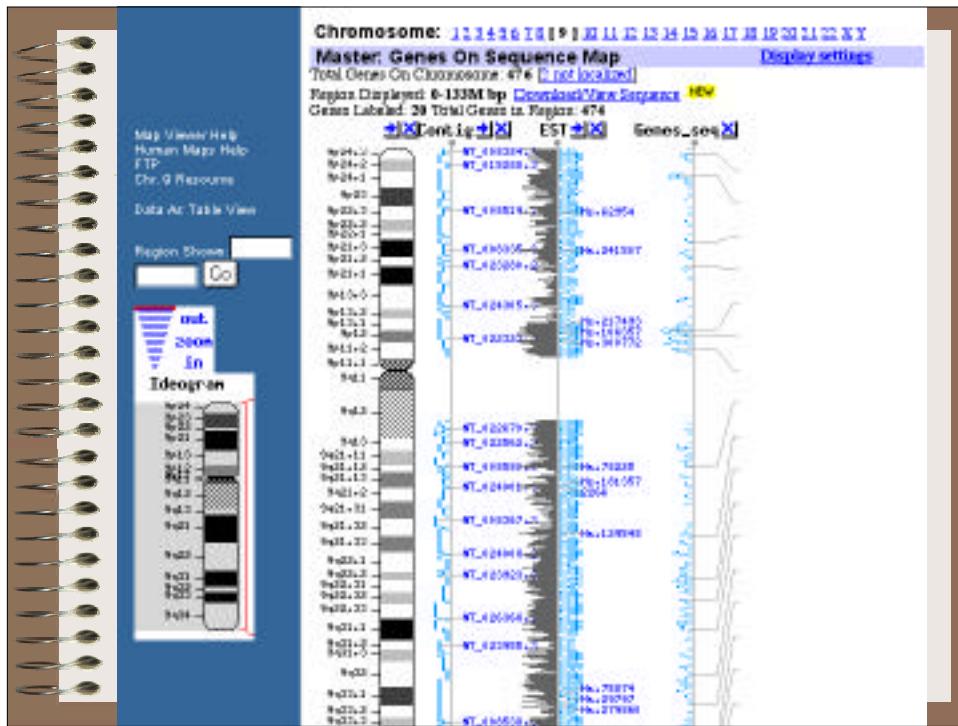
When browsing the genome using the new MapViewer, click on Display Settings to choose from several types of maps and . Below are three views of the BRCA2 locus using different display options. Click the image to see the full MapViewer display.

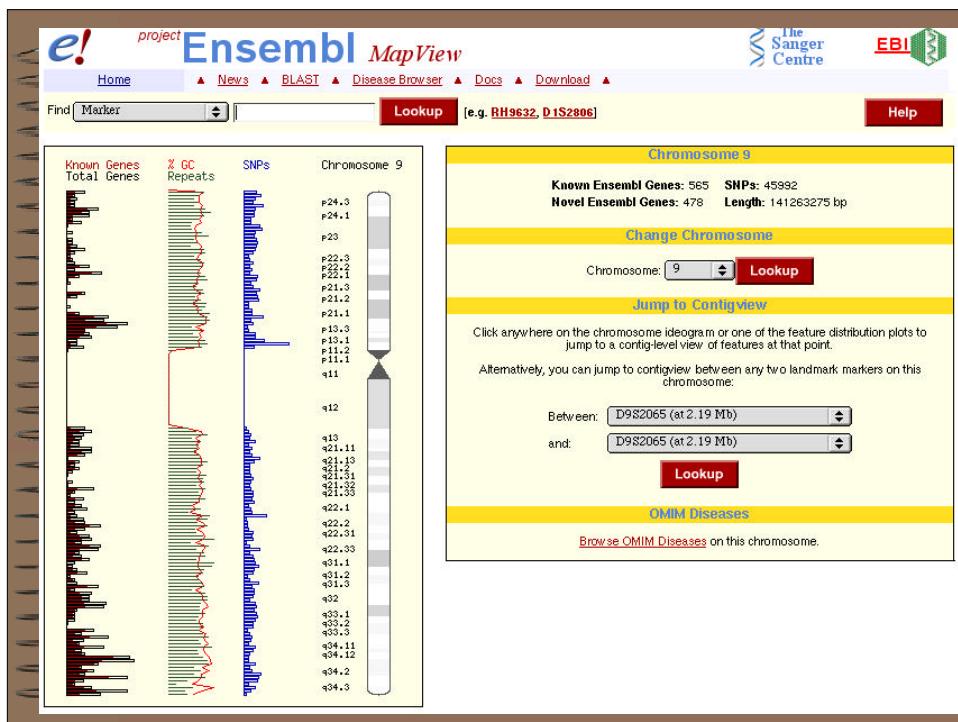
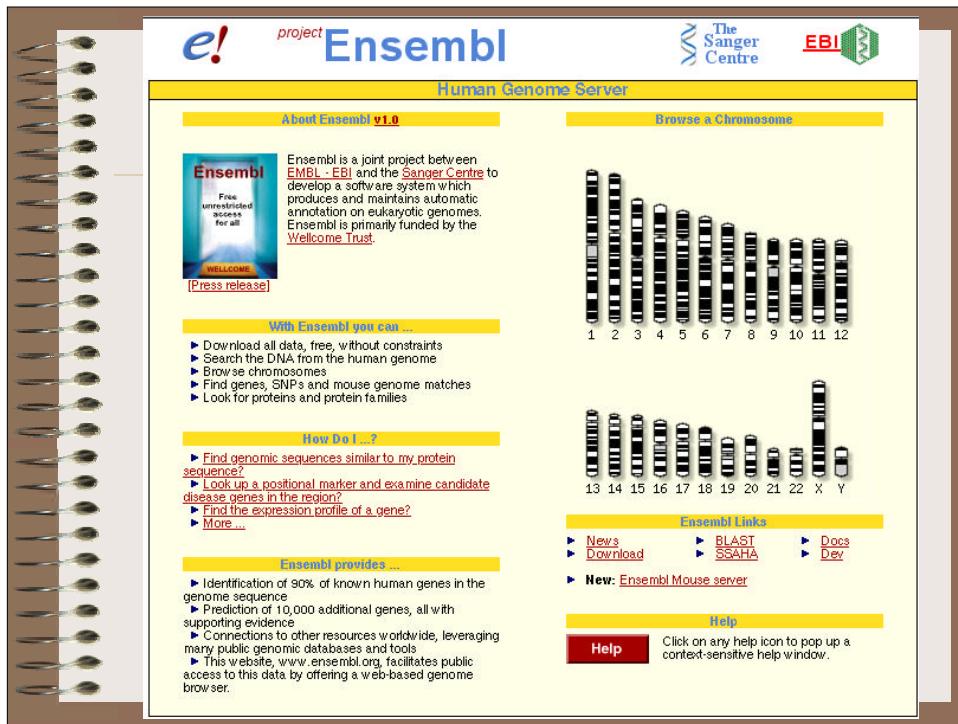
NCBI Display options

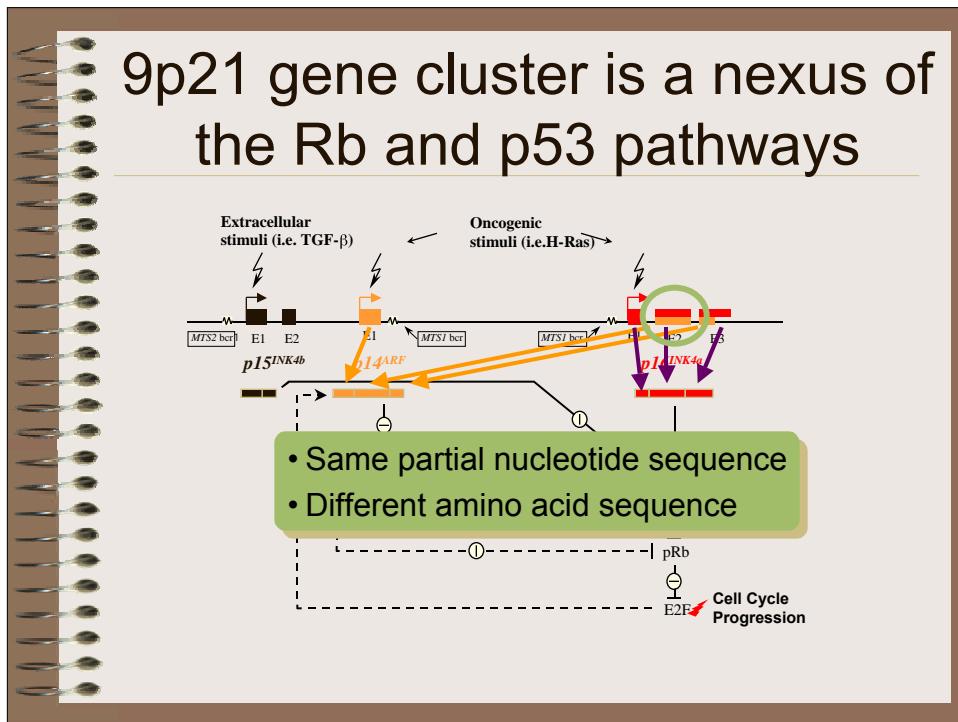
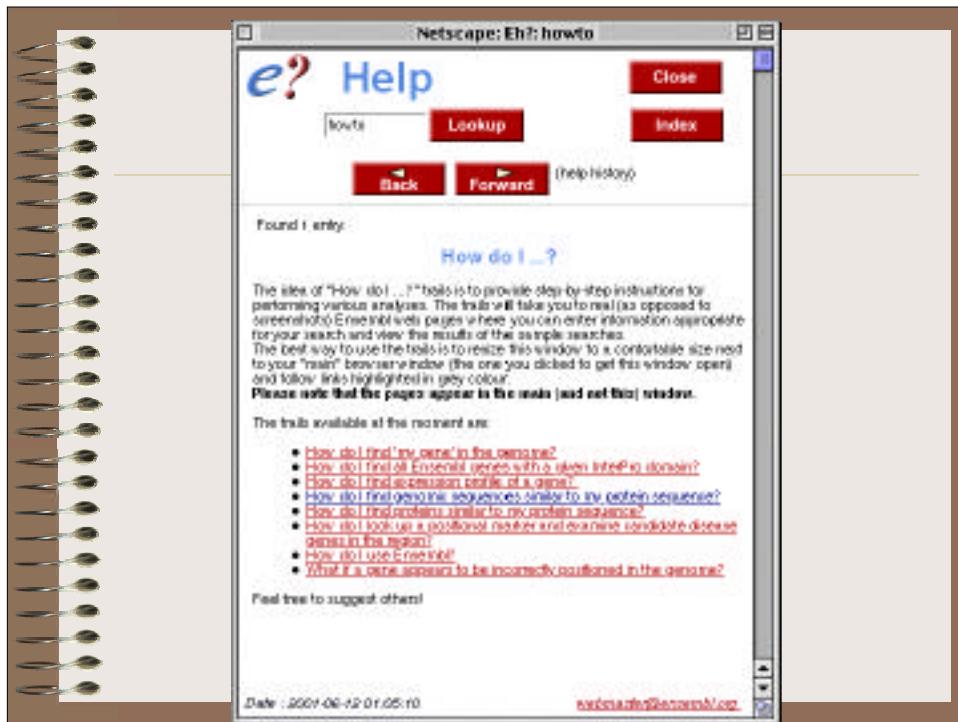
A. Genes

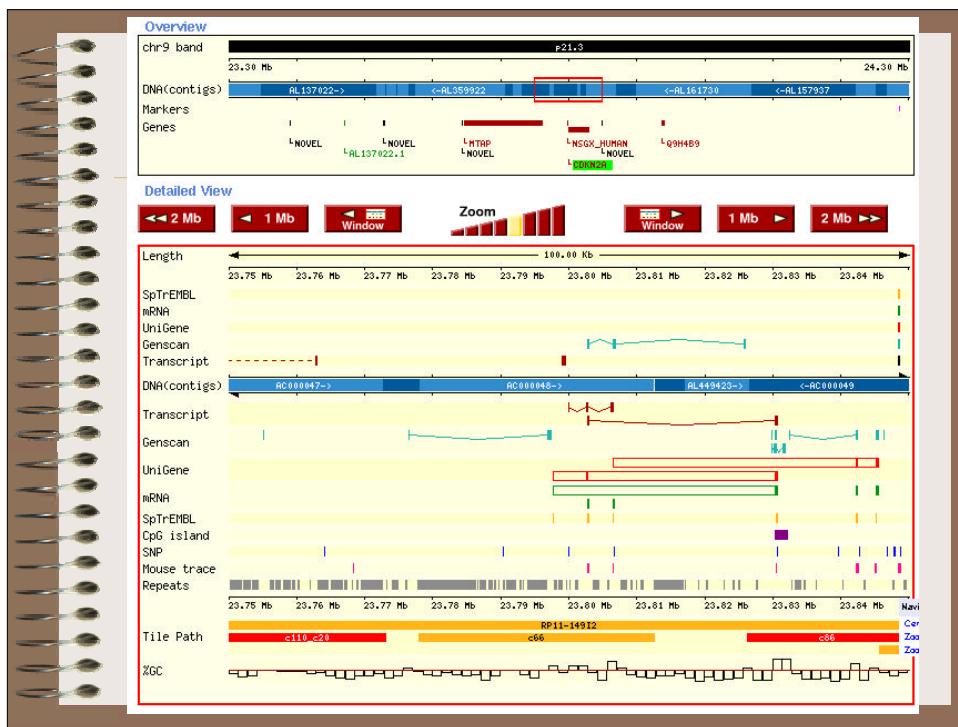
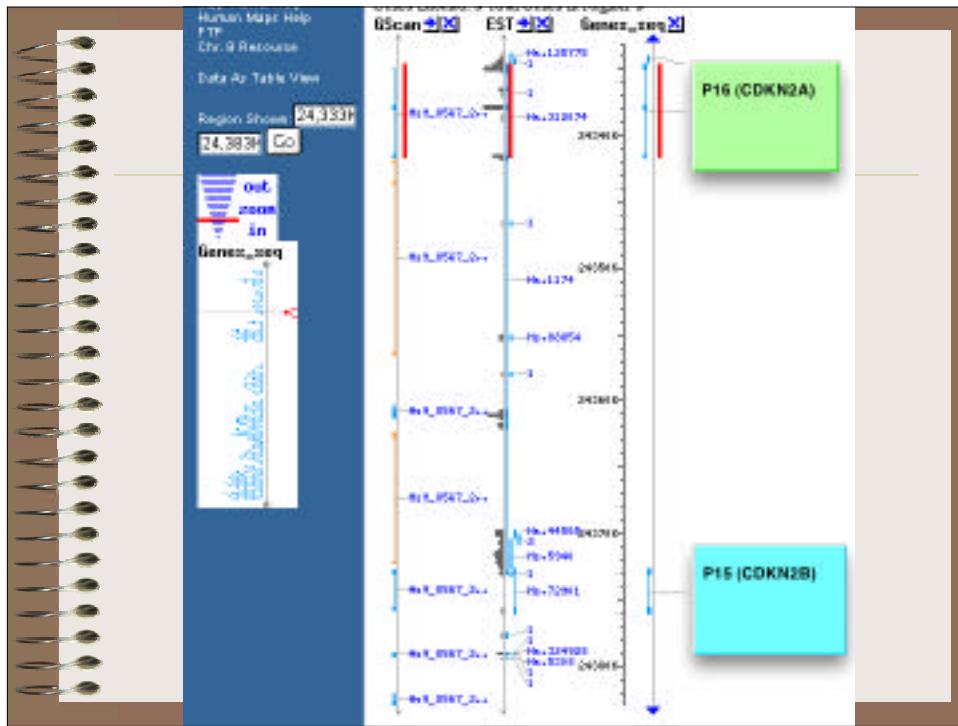
B. Variations, juxtaposed with genes

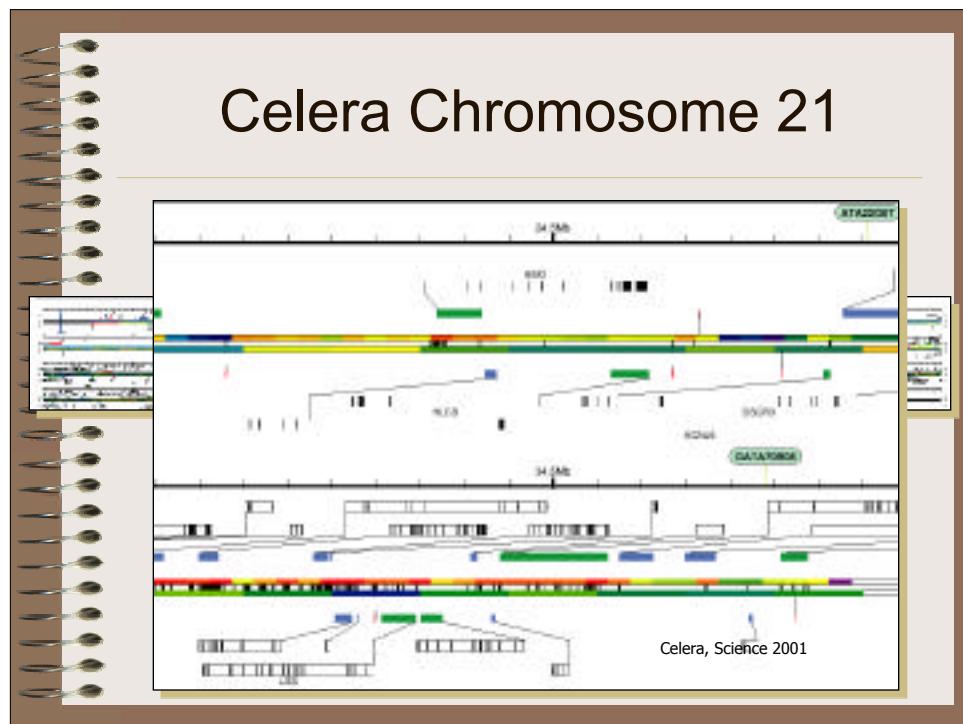
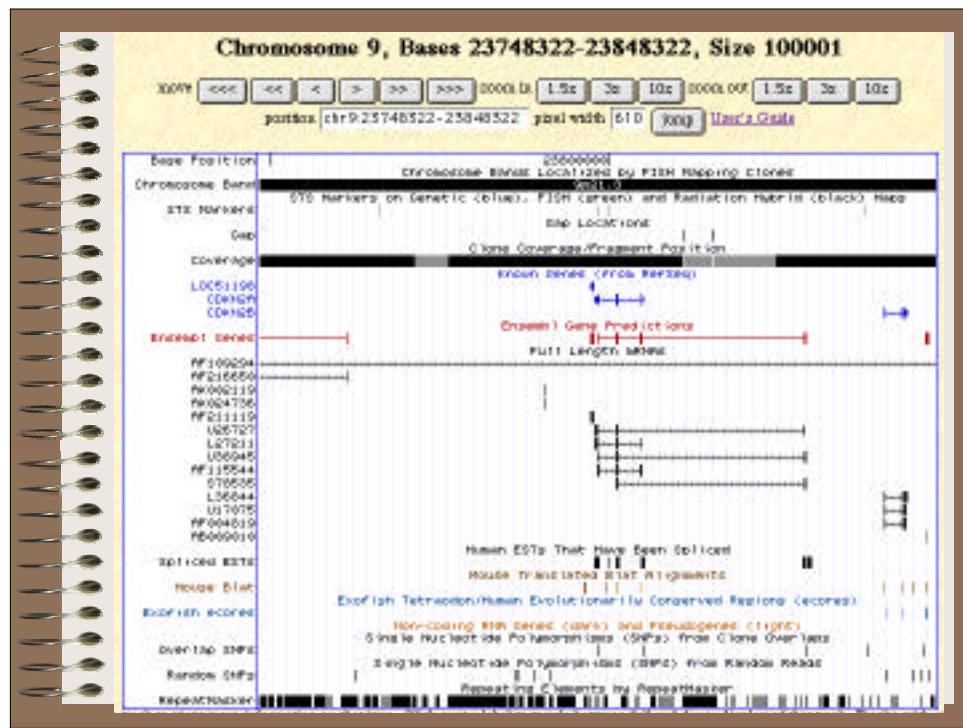
C. Several STS maps, juxtaposed with genes

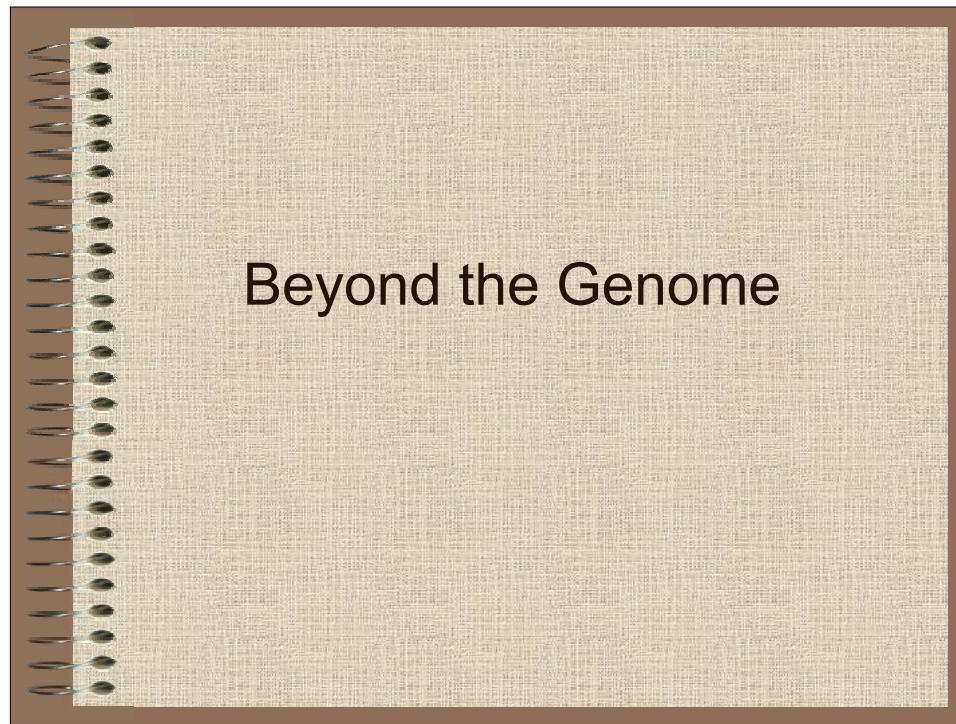












 **ExPASy Molecular Biology Server**

This is the ExPASy (Expert Protein Analysis System) proteomics server of the [Swiss Institute of Bioinformatics](#) (SIB). This server is dedicated to the analysis of protein sequences and structures as well as 2-D PAGE ([Disclaimer](#)).

[[Announcements](#)] [[Job openings](#)] [[Mirror Sites](#)]

| Databases | Tools and Software Packages |
|---|---|
| <ul style="list-style-type: none"> • SWISS-PROT and TrEMBL - Protein sequences • PROSITE - Protein families and domains • SWISS-2DPAGE - Two-dimensional polyacrylamide gel electrophoresis macromolecules • SWISS-3DIMAGE - 3D images of proteins and other biological macromolecules • SWISS-MODEL Repository - Automatically generated protein models • CD40Lbase - CD40 ligand defects • ENZYME - Enzyme nomenclature • SeqAnalRef - Sequence analysis bibliography references • Links to many other molecular biology databases | <ul style="list-style-type: none"> • Proteomics tools <ul style="list-style-type: none"> ◦ Identification and characterization ◦ DNA -> Protein ◦ Similarity searches ◦ Pattern and profile searches ◦ Post-translational modification prediction ◦ Primary structure analysis ◦ Secondary structure prediction ◦ Tertiary structure ◦ Transmembrane regions detection ◦ Alignment • Molano 3 - Software for 2-D PAGE analysis • SWISS-MODEL - Automated knowledge-based protein modeling server • Swiss-PubViewer - Macintosh/PC tool for structure display and analysis • Boehringer Mannheim's Biochemical Pathways |
| Education and services <ul style="list-style-type: none"> • The ExPASy FTP server • Swiss-Shop - automatically obtain (by email) new sequence entries relevant to your field(s) of interest • Masters Degree in Bioinformatics • 2-D PAGE training - attend a one-week course in Geneva • SWISS-2DSERVICE - get your 2-D Gels performed according to Swiss standards | Documentation <ul style="list-style-type: none"> • What's New on ExPASy • SWISS-FLASH electronic bulletins • SWISS-PROT documents • How to create HTML links to ExPASy • Complete table of available documents |
| Links to lists of molecular biology resources <ul style="list-style-type: none"> • Amos' WWW links - The ExPASy list of Biomolecular servers • BioHunt - Search the internet for molecular biology information • WORLD-2DPAGE - Links to 2-D PAGE database servers and 2-D PAGE related servers and services | Links to some major molecular biology servers <ul style="list-style-type: none"> • European Bioinformatics Institute (EBI) • National Center for Biotechnology Information (NCBI) • Japanese GenomeNet • Australian National Genomic Information Service |

Physical Properties

Prediction of Physical Properties

- Compute pI/MW <http://www.expasy.ch/tools/pitool.html>
- MOWSE <http://srs.hgmp.mrc.ac.uk/cgi-bin/mowse>
- PeptideMass <http://www.expasy.ch/tools/peptide-mass.html>
- TGREASE <ftp://ftp.virginia.edu/pub/fasta/>
- SAPS <http://www.isrec.isb-sib.ch/software/SAPSform.html>

Prediction of Protein Identity Based on Composition

- AACompIdent <http://www.expasy.ch/tools/aacomp/>
- AACompSim <http://www.expasy.ch/tools/aacsim/>
- PROPSERACH <http://www.embl-heidelberg.de/prs.html>

Motifs and Patterns

- BLOCKS <http://blocks.fhcrc.org>
- Pfam <http://www.sanger.ac.uk/Software/Pfam/>
- PRINTS <http://www.bioinf.man.ac.uk/dbrowser/PRINTS/PRINTS.html>
- ProfileScan <http://www.isrec.isb-sib.ch/software/PFSCANform.html>

Protein Structure

Predictor of Secondary Structure and Folding Classes

- nnpredict http://www.cmpharm.ucsf.edu/_nomi/nnpredict.html
- PredictProtein <http://www.embl-heidelberg.de/predictprotein/>
- SOPMA <http://pbil.ibcp.fr/>
- Jpred <http://jura.ebi.ac.uk:8888/>
- PSIPRED <http://insulin.brunel.ac.uk/psipred>
- PREDATOR <http://www.embl-heidelberg.de/predator/predatorinfo.html>

Prediction of Specialized Structures or Features

- COILS <http://www.ch.embnet.org/software/COILSform.html>
- MacStripe <http://www.york.ac.uk/depts/biol/units/coils/mstr2.html>
- PHDtopology <http://www.embl-heidelberg.de/predictprotein>
- SignalP <http://www.cbs.dtu.dk/services/SignalP/>
- TMpred <http://www.isrec.isb-sib.ch/ftp-erver/tmpred/www/TMPREDform.html>

Structure Prediction

- DALI <http://www2.ebi.ac.uk/dali/>
- Bryant-Lawrence <ftp://ncbi.nlm.nih.gov/pub/pkb/>
- FSSP <http://www2.ebi.ac.uk/dali/fssp/>
- UCLA-DOE <http://fold.doe-mbi.ucla.edu/Home>
- SWISS-MODEL <http://www.expasy.ch/swissmod/SWISS-MODEL.html>
- TOPITS <http://www.embl-heidelberg.de/predictprotein/>